



## MOve: Design of an Application-Malleable Overlay

Sébastien Monnet, Ramsés Morales, Gabriel Antoniu, Indranil Gupta

### ► To cite this version:

Sébastien Monnet, Ramsés Morales, Gabriel Antoniu, Indranil Gupta. MOve: Design of an Application-Malleable Overlay. [Research Report] RR-5872, INRIA. 2006, pp.18. inria-00070154

**HAL Id: inria-00070154**

**<https://hal.inria.fr/inria-00070154>**

Submitted on 19 May 2006

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET EN AUTOMATIQUE

## ***MOve: Design of an Application-Malleable Overlay***

Sébastien Monnet , Ramsés Morales , Gabriel Antoniu and Indranil Gupta

**N°5872**

Mars 2006

————— Systèmes numériques —————

A large blue rectangle occupies the lower half of the page. Overlaid on it is a large, light gray stylized 'R' logo. To the right of the 'R', the words 'Rapport de recherche' are written in a white serif font. A horizontal gray brushstroke is positioned below the text.

***Rapport  
de recherche***





## MOve: Design of an Application-Malleable Overlay

Sébastien Monnet <sup>\*</sup>, Ramsés Morales <sup>†</sup>, Gabriel Antoniu <sup>‡</sup> and Indranil Gupta <sup>§</sup>

Systèmes numériques  
Projet Paris

Rapport de recherche n°5872 — Mars 2006 — 18 pages

**Abstract:** Peer-to-peer overlays allow distributed applications to work in a wide-area, scalable, and fault-tolerant manner. However, most structured and unstructured overlays present in literature today are *inflexible* from the application viewpoint. In other words, the application has no control over the structure of the overlay itself. This paper proposes the concept of an *application-malleable* overlay, and the design of the first malleable overlay which we call *MOve*. In *MOve*, the communication characteristics of the distributed application using the overlay can *influence* the overlay's structure itself, with the twin goals of (1) optimizing the application performance by adapting the overlay, while also (2) retaining the scale and fault-tolerance of the overlay approach. The influence could either be explicitly specified by the application or implicitly gleaned by our algorithms. Besides neighbor list membership management, *MOve* also contains algorithms for resource discovery, update propagation, and churn-resistance. The emergent behavior of the implicit mechanisms used in *MOve* manifest in the following way: when application communication is low, most overlay links keep their default configuration; however, as application communication characteristics become more evident, the overlay *gracefully* adapts itself to the application.

**Key-words:** Peer-to-peer, overlay, malleable, gossip

(Résumé : tsvp)

This work was supported by NSF CAREER grant CNS-0448246 and UIUC-INRIA collaboration grant.

<sup>\*</sup> Sebastien.Monnet@irisa.fr

<sup>†</sup> rvmorale@cs.uiuc.edu

<sup>‡</sup> Gabriel.Antoniu@irisa.fr

<sup>§</sup> indy@cs.uiuc.edu

## **MOve: conception d'un réseau logique malléable**

**Résumé :** Les réseaux logiques pair-à-pair permettent aux applications distribuées de s'exécuter à grande échelle de manière complètement distribuée et tolérante aux fautes. Cependant, la plupart des réseaux logiques présentés dans la littérature, structurés et non-structurés, sont inflexibles du point de vue de l'application. En d'autres termes, l'application n'a pas de contrôle sur la structure du réseau logique qu'elle utilise. Ce papier propose le concept d'un réseau logique malléable, et la conception du premier réseau logique malléable que nous appelons "MOve" (pour Malleable Overlay). Les caractéristiques des patterns de communication d'une application distribuée utilisant MOve, peuvent influencer la structure même du réseau logique. Ceci avec un double but: 1) optimiser les performances des applications en adaptant le réseau logique, tout en 2) conservant les propriétés de passage à l'échelle et de tolérance aux fautes des réseaux logiques. Cette influence peut être spécifiée explicitement par l'application, ou glanée par nos algorithmes. En plus de gérer des listes de voisins (constituant le réseau), MOve contient également des algorithmes de découvertes de ressources, de propagation des mises à jour et de résistance à la volatilité. Quand l'application a peu de besoins de communications, la plupart des liens du réseau logique conservent leur configuration par défaut; cependant, quand les schémas des communications de l'application deviennent plus évidents, le réseau logique s'adapte de lui-même en favorisant les liens entre les entités communicantes.

**Mots-clé :** Pair-à-pair, réseaux logiques, malléable

## 1 Introduction

Today, peer-to-peer (P2P) overlays fall into two categories - (1) structured (i.e. *Distributed Hash Table*-based) overlays such as Pastry and Chord [12, 14], and (2) unstructured (i.e., gossip- or flooding-based) overlays such as Freenet, Gnutella, KaZaA [5, 11, 18]. These P2P overlays offer reliability in the face of massive failures and churn (node join and leave), as well as scalability to hundreds and thousands of nodes.

However, both these types of overlays have the common disadvantage that they are *inflexible* from the application viewpoint. The rules and invariants for selecting and maintaining neighbor nodes in the overlay, as well as resource discovery, are all dictated in a rigid fashion by the underlying overlay. This usually means that the developer of a distributed application has a limited number of options – either go with the provided overlay, or design a new overlay from scratch.

In this paper, we propose the concept of an *application-malleable* overlay. An application-malleable overlay is defined as an overlay where the communication characteristics of the distributed application using the overlay can *influence* the overlay’s structure itself. The twin goals of a malleable overlay are: (1) optimizing application performance from the overlay, while also (2) retaining the scale and fault-tolerance of the overlay approach.

In order to realize and evaluate our design philosophy, we build a specific malleable overlay called *MOve* (for *Malleable OVERlay*) that combines elements of an *unstructured overlay* with application characteristics. In *MOve*, the structure and behavior of the overlay is influenced both by the underlying default unstructured overlay, as well as *influenced* by application characteristics. The influence could either be (1) explicitly specified by the application or (2) implicitly gleaned by our algorithms.

In a P2P overlay, each node maintains a separate neighbor list - this is a membership list that specifies the *who knows whom* relationship. This neighbor list is partial in the sense that it contains only some of the nodes in the system [6, 4, 14]. *MOve* contains algorithms for neighbor list maintenance, propagating updates, and churn-resistance. The most interesting feature of *MOve* is its emergent behavior. When application communication is low, *MOve* autonomically evolves to keep most of the overlay links in the default state so that most of the system looks like an unstructured overlay. However, as application communication characteristics become more and more evident, the overlay autonomically *gracefully* adapts itself to the application, but without forgetting its default structure.

To focus our approach on a particular class of applications, we choose *collaborative applications*, such as distributed whiteboard platform, an audio/video conferencing service, a replicated data-sharing service, or a distributed-gaming platform. All these applications rely on the notion of *application groups* - each process belongs to one or more groups, and interacts with other processes in common groups. For instance, the members of the same group may share a distributed state that needs to be updated (the whiteboard, the game-board, the replicas of a mutable piece of data, etc.). Alternately, the set of replica managers for a particular data item would form a group.

*MOve* allows such *group*-based applications to influence the underlying overlay, that may be common to multiple, coexisting, collaborative applications. Besides the neighbor list maintenance, *MOve* has the following three goals: (1) (*connectivity*) maintaining a low diameter for the overlay so that unstructured queries can be propagated quickly; (2) (*updates*) efficient update propagation

mechanisms within and across groups; and (3) (*volatility-resilience*) combat volatility arising from rapid node arrival and failure (i.e., “churn”).

The basic idea in the MOve approach is to have each node maintain a neighbor list that, by default, consists of *non-application neighbors*, i.e., randomly selected neighbors. However, with the formation of more and more application groups, some of these non-application neighbors are automatically replaced by *application-aware* neighbors (shortly: *application neighbors*). A non-application neighbor may either turn into an application neighbor (if the neighbor belongs to a common group) or be replaced by a new application neighbor. We have implemented MOve, and our experiments show (1) that the system manages to maintain logarithmic overlay path lengths; (2) that it gracefully manages transitions between application and non-application neighbors as the number of groups increases and decreases. In addition, MOve shows good scalability and volatility-resilience.

The next section presents research efforts related to the various aspects involved in this context. Then Section 3 describes a scenario explaining why it is important to have efficient communication within groups of nodes. Section 4 gives a general overview of our approach, and provides an analysis. Section 5 presents some preliminary simulation results. Finally, Section 6 discusses the contribution and the future work.

## 2 Related Work

In the past few years, many research efforts have focused on building overlays for peer-to-peer networks, essentially for large-scale immutable file-sharing. For this kind of application, predicting which node is going to communicate with which node is not trivial. Therefore, most algorithms for building overlays do not take communication patterns into account.

In unstructured (i.e., gossip- or flooding-based) P2P overlays, such as [11] or [15], neighbor lists are usually built and maintained by randomly selecting a subset of the neighbors’ neighbors. The goal is therefore to obtain a random graph.

In structured (i.e., DHT-based) overlays, the *who knows who* relation is usually defined by means of a given topology (typically a ring); the position of each node in this given topology is determined by a hash function on its IP address [12, 14].

Even if some of these previous proposals take into account certain criteria while building the overlay, (e.g., physical locality in the case of [12]), they do not take into account the application-related relations between nodes, which can express interaction patterns that may result from the way the overlay is solicited by the application.

Very few recent research efforts take into account these relations. Semantic overlay networks [7, 16] exploit the semantic relations between peers (based on the set of files they share). They propose solutions allowing to improve the efficiency of the search mechanisms for large-scale file-sharing applications, by creating shortcuts between peers which are semantically close. Efficient search is also the goal addressed by the *path-caching* technique, which consists in keeping data references along a given search path, in order to improve the efficiency of subsequent search operations.

However, for the group-based applications we target in this paper, such as distributed shared whiteboard (or gaming) platforms, or replicated data-sharing services, search efficiency is not the only property to optimize. In such applications, the members of a same group share some data (a whiteboard, a replicated piece of data, a game state, etc.), which they all potentially read and write.

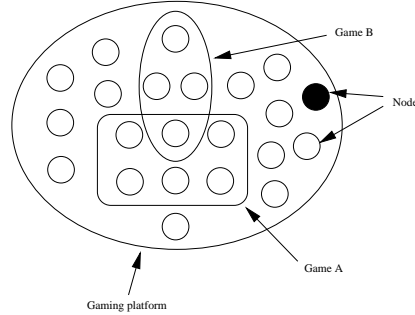


Figure 1: A gaming platform

When a peer writes this shared data, it is important for the updates to be efficiently propagated to the other members of the group. Consequently, the peers belonging to a same group have to be close to each other in the overlay (i.e., a few hops away), to enable the application to efficiently maintain the consistency of the members' views of the shared state.

The issue of update propagation in large-scale systems has been studied in [13]. This system proposes an efficient multicast scheme based on multicast trees built on top of the Pastry overlay. The problem we address in this paper is different, since every member of the group can be the source of multicast in our target applications. The approach we propose is also different, since it does not construct a membership mechanism based on an already *existing overlay*. Our goal is to build an *emergent and adaptive overlay* based on patterns derived from the application usage.

The closest work related to ours is [9], which addresses the problem of building an adaptive overlay based on different criteria: topology, semantic proximity, bandwidth, etc. The problem is addressed in a generic way: the target scale and the target applications are not specified. The issue we address is more specific: it regards applications that need efficient updates within groups of nodes. Consequently, multiple criteria have to simultaneously be taken into account and controlled: application-dependent node relations, but also physical locality, as well as the connectivity of the resulting graph (expressed through the degree of clustering).

### 3 Scenario

To motivate our work, we consider a large scale distributed gaming platform (represented by Figure 1). This application may involve tens of thousands of nodes spread around the Internet. For efficiency reasons, the number of neighbors that a peer must know has to be bounded, since the related information requires monitoring and state updating. Therefore, each node only has a partial view of the system. However, this should not have a negative impact on the application's desired properties, such as *connectivity*, *update efficiency* and *volatility resilience*, which are important for collaborative applications, such as gaming platforms.

**Connectivity.** Some particular node (for instance the black node in Figure 1) may have to lookup for a specific game instance in the platform (e.g., Game A). This game may involve only a small subset of nodes (few tens). The neighbors contained in this particular node's neighbor list may not



be involved in this particular game. The platform has to be connected to make it possible for a node to reach somehow (even through a quite long path) all the other nodes. The lookup of a node that is participating in Game A is application dependent: it could be done by visiting a website, or querying a custom search engine, or by flooding a search query on the overlay.

For efficiency purposes, the diameter of the overlay should be as small as possible, even with partial neighbor lists. To achieve this, the graph formed by the nodes and links needs to have well distributed degrees. Note that it is enough for a new player to find only one player for the wished game in order to be able to reach the other ones.

**Efficient Updates.** While a game is running, the players store object replicas which represents the current state of the game (depending on the application, this can correspond to a shared white board, a piece of replicated data, etc). Each time a player plays, his node updates the state of its local game board version (i.e., its replica). In order for the other players to be able to play, they have to be notified of the changes in the game board. Therefore, updates need to be propagated in an efficient manner within a group (e.g., Game A or Game B in Figure 1).

**Volatility Resilience.** Among several thousands of nodes spread over the Internet, it is likely that from time to time some nodes fail or get disconnected. At the global level (the entire platform), failures and disconnections may not lead to break down the whole graph connectivity.

At the level of a given game (i.e., groups) such events should not stop the game, which means that the remaining players have to remain connected together. Furthermore, the departure of one player may break some path in the group graph. The longest path between two nodes (the subgraph diameter) is likely to grow; however, the update propagation mechanism has to stay efficient.

## 4 Design

The first purpose of an overlay is to connect nodes together. Therefore, the first property to fulfill is the *connectivity* of the constructed graph. Furthermore, Section 3 highlighted the importance to provide the ability to perform efficient updates among groups of nodes within the overlay. This implies that some *clustering* is necessary. Both *connectivity* and *clustering* have to be preserved while taking into account the dynamic nature of the environment.

**Random Graph Benefits.** Graph theory shows that *random graphs* have good properties in terms of *connectivity* and *degree distribution*. For instance, in a random graph, if each node has at least  $\log(N)$  uniformly random neighbors (where  $N$  is the total number of nodes) the random graph will be connected with high probability [2]. Estimating the size (i.e.,  $N$ ) of a large scale dynamic distributed system has also been addressed in previous studies [10]. In our case, the scale is not infinite (we target thousands to a few tens of thousands of nodes), therefore safe bounds can be assumed instead. For instance, 50 links per node will provide a large safety margin to theoretically connect  $5 * 10^{21}$  nodes. On the other hand, random graphs also have the benefit of leading to a good degree distribution. An overlay based on a random graph may take advantage of this for *load distribution*.

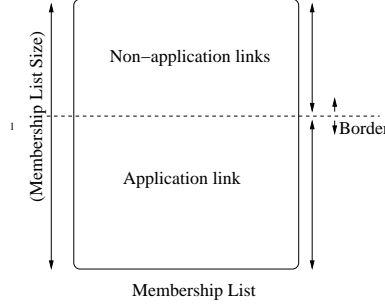


Figure 2: The neighbor list on each node

In our design, nodes maintain a neighbor list, containing links to the node's neighbors. For each node, an upper bound ( $l$ ) is set on the size of the neighbor list. This bound is first set according to an initial approximation of the network size (while observing the condition  $l > \log(N)$ ). Then, during the execution, this value can be increased when necessary, if allowed by the available resources (see below).

The neighbor list is composed of two kinds of links: *non-application links* and *application links*. Figure 2 represents a node's neighbor list.

**Non-Application Links.** *Non-application links* are responsible for maintaining a global overlay, with a low degree of clustering. If the application is in a state that does not need clustering (e.g., at initialization), the neighbor list will contain only non-application links. Remember that nodes don't need full knowledge about the network, and the number of non-application links may vary from node to node.

**Application Links.** To cluster together nodes that belong to a group  $i$ , each member of the group creates  $k_i$  *application links* to other randomly chosen members of the same group. This clustering will enable fast propagation of state updates among the members of the group. It will also favor an efficient propagation of application-level multicast messages. Parameter  $k_i$  is determined by the application, and it must be at least  $\lceil \ln(|R_i|) \rceil$ , where  $R_i$  is the number of members of group  $i$ . Essentially, the goal is to create a strongly connected graph for group  $i$  (with a short characteristic path length).

**Replacement Policy.** When an application link needs to be created, it will be added to the neighbor list following four different ways. Assume that we want to create an application link for group  $i$ , that points to node  $n$ . (1) If the size of the neighbor list is smaller than  $l$ , and there is no non-application link pointing to  $n$ , a new link will be added to the list. (2) If the size of the neighbor list has already reached  $l$ , but the node has enough resources available, the neighbor list can grow to accommodate the link (by setting  $l$  to a larger value for the current node). (3) If the node decides not to grow the list, then a non-application link will be dropped, and the application link will be added. (4) Finally, if there is a non-application link pointing to  $n$ , then it will become an application link.

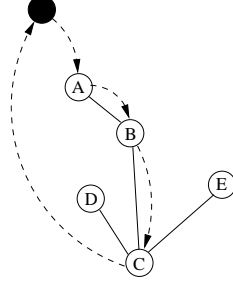


Figure 3: Join mechanism

#### 4.1 Addressing the Connectivity Issue

An overlay is said to be *connected* if there is a path (succession of edges or links) between every pair of nodes. This property is very important in an overlay as it provides the guarantee for a node to be able to communicate with all the other ones in the overlay.

The tradeoff between the good properties of random graphs and those enabled by favoring clustering between related nodes can be tuned by setting some bounds. The first bound is the size  $l$  of the neighbor list, which is managed as explained above. The second bound is  $k_i$ , which limits the number of links that are involved in the group  $i$ . If  $l - \sum_i k_i$  is large enough (a few tens for the scale we are targeting), i.e., the neighbor list contains enough non-application links, the good properties of random graphs are available in spite of the little clustering induced by taking the topology into account. On the other hand, it is important to notice, that  $\sum_i k_i$  is in fact greater than the number of *application links*. This is explained by the fact that some *application links* may be shared by multiple groups when there are group intersections.

**A node joins** the overlay (e.g., the black node in Figure 3) by contacting any current overlay member. If the peer that receives the join request has space available in its neighbor list, it will reply with its current neighbor list, and will add a link to the joining node to its non-application links. If the neighbor list is full (which is the case for Nodes A and B on Figure 3), the join request will be forwarded to a randomly chosen node. The forwarding of a join request is associated with a time to live (TTL). If all nodes that receive the forwarded request have full neighbor lists, the TTL will reach 0, and the last node to receive the forward will forcibly add a link to the new node to its non-application links. It will then reply with its current neighbor list (e.g., Node C). We do this to ensure that the in-degree of a node is always above 0. The new node will use the received neighbor list to create its own list.

The **failure detection protocol** is based on the SWIM [6] protocol. Each protocol period, of length  $T$  seconds<sup>1</sup>, each node sends a ping message to one of its neighbors. The target node is selected by sequentially traversing an array that represents the random permutation of the neighbor list. Once the array is completely traversed, a new permutation is computed. The node expects a reply to the ping message within a timeout of  $t < T$  seconds. If the reply is not received on time, an indirect ping is sent to  $y$  nodes. These nodes will then send a ping to the intended target node, and,

<sup>1</sup>Protocol periods are asynchronous at different process, although it is assumed that they have the same  $T$ .

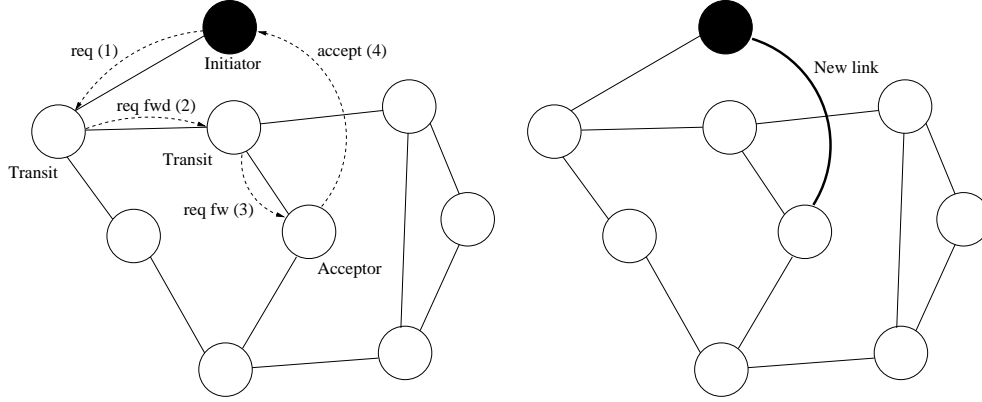


Figure 4: The random-walk mechanism

if they receive a reply, the reply will be sent back to the node that originated the ping. The intention of the indirect ping is to sidestep transient network problems. If no reply is received before the next protocol period, the ping target will be *suspected* of having failed. At the beginning of each protocol period, any node that has been suspect for one protocol period will be dropped from the neighbor list, i.e., declared dead.

To achieve an overlay with a low clustering coefficient and evenly distributed in-degree, every  $U$  protocol periods, each node verifies if its non-application membership list has been modified. If no modification has been made after  $U$  periods, it issues a join message to a random node. With the membership list it will receive as a reply to its join message, the node will try to replace a fraction<sup>2</sup> of its own membership list. Note that the smaller  $U$  is, the more aggressive the replacement will be, and the faster the protocol will take the overlay to a stable low-clustering coefficient.

## 4.2 Group Communications

As previously explained, when a new application link is created, it will result in the substitution of a non-application link, unless the node has enough resources to grow its neighbor list. In this way we maintain the overhead constant at the node. On the other hand, an application link can be *shared*. For instance, assume Node  $a$  belongs to Groups  $i$  and  $j$ . If Node  $b$  joins Groups  $i$  and  $j$ , it creates a single application link to  $a$ , knowing that this link is *shared*. A sharing count is maintained for such links.

**Random Walk for Application Links.** To cope with the dynamic nature of the infrastructure and avoid pathological topologies that may be induced by failures, it is important to periodically refresh the links. This is also useful in order to guarantee a small path between any two nodes in a given group. To this effect, we rely on another result from random graph theory [2]: adding  $O(n)$  non-application links to a graph with  $n$  vertices will reduce the diameter to  $O(\log(n))$ . This result only

<sup>2</sup> $f = 50\%$  in our evaluations

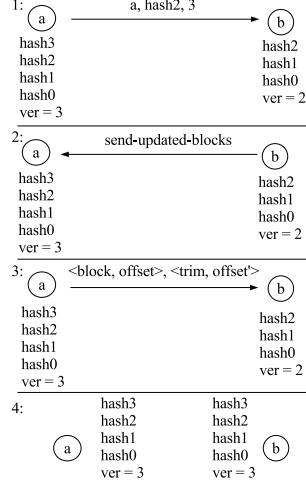


Figure 5: 1: *a* sends an update message to *b*. *a* contains 4 hashes on its hash list, and *b* contains 3 hashes on its hash list that coincide with *a*'s three oldest hashes. 2: the hash included in the message coincides with the newest hash at *b*, so *b* requests an updated blocks transfer. 3: *a* sends the updated blocks, and in this example we also show deletion of the tail of the replica (*trim* at offset *offset'*). 4: after the update is applied, both nodes have the same hash list and version numbers.

applies to undirected graphs. Thus we add the restriction that all application links are *bidirectional*. When an application link is created from *a* to *b*, *b* will also create a link to *a*. If node *b* deletes the link, so will node *a*. When an application link is shared, it will be maintained until the sharing count reaches 0. When an application link stops being used as such, it is changed to a non-application link. This simulates an undirected graph inside the application group.

The graph is refreshed periodically, by having every node in an application group execute the following steps:

- (1) Launch a random walk to get a new neighbor. The random walk hops at most *TTL* times, using application links that belong to the group.
- (2) Drop an old link when the new link is created.

Although it is assumed that the timeout to launch the random walk is an application parameter, note that nodes that belong to a group are not synchronized. Also note that the bidirectionality of the links is always enforced.

**Update Propagation.** The update propagation algorithm is designed to propagate and apply updates under a dynamic set of nodes and object replicas, and it provides causal consistency [8]. The update is transferred across members of the group. Although the update mechanism detects update conflicts, the actual resolution is deferred to the application layer, which receives an alarm with the conflict.

When a node updates a replica, a secure hash is computed over the entire replica contents, and an object version number is incremented. The hash is added to a list of hashes of previous replica versions –this list is ordered chronologically. The node then broadcasts an update message using the *application links* associated with the group. The message contains the IP number, the new version number, and the replica hash that corresponds to the replica before the update. When a node receives the message it verifies the version number, if it is lower than the local version number then the message is ignored. If the message is not dropped, it is forwarded on all outgoing *application links* that correspond to the group. After forwarding the update message, the node compares its current replica hash with the one received on the message, if they are equal then the node will request the upstream node to forward the updated blocks of the replica. If the hash differs the node will send its current replica hash to the upstream node. The upstream node will search for this hash in its hash list, and if found, the upstream node will copy the entire replica to the downstream node along with the list of hashes, otherwise a conflict alarm will be sent to the application layer. Figure 4.2 shows a common case example.<sup>3</sup>

The application can choose a lazy update propagation approach. In this case the update messages will be forwarded as before, but the actual update will not be propagated until a node needs to write or read its replica. When a node reads or writes, it will ask the upstream node for the updated data, which could itself request the update to another upstream node.

### 4.3 Analysis

To show the benefits of link-sharing among the different application groups and non-application links, we present an analysis of a variant of the MOve system. This variant does not impose limits on neighbor list sizes, and allows them to grow indefinitely. *Without* the MOve approach, this indefinite approach would have the neighbor list size of each node grow linearly as the sum of the number of its neighbors for each application group the node belongs to, and the number of its non-application neighbors. With our MOve approach, the number of links saved are significant, as shown by our analysis below.

Formally, at a node  $p$ , let  $Nbrs_N(p)$  represent the set of application links at  $p$ . Assume that  $p$  belongs to  $k$  groups  $R_i$  ( $i = 1$  to  $k$ ). Let  $Nbrs_{|R_i|}(p)$  represent the set of neighbors that  $p$  has in group  $R_i$ . Now, let  $f(N) = |Nbrs_N(p)|$  and for each  $i = 1$  to  $k$ , let  $f(R_i) = |Nbrs_{|R_i|}(p)|$ . Note that in our implementation,  $f(x) = O(\log(x))$ , however our results are more general.

**Theorem 1:** Assume that: (1) each application group consists of members selected uniformly at random, and (2) at node  $p$ , each non-application link for a group is selected uniformly at random among the group members and (3) at node  $p$ , each application link for a group is selected uniformly at random from among the group members<sup>4</sup> Then: (a) the expected number of links saved by MOve is positive and (b) it grows linearly as the number of non-application links is increased, and (c) it is

<sup>3</sup>We do not implement version vectors on replicas, to separate real conflicts from updates that can be merged, because that approach would not allow our system to scale as intended.

<sup>4</sup>The uniformly at random assumption number (3) is reasonable since our neighbor list maintenance protocols achieve such random neighbor lists.

proportional to the total number of non-application links if all groups at node  $p$  are equi-sized.

**Proof:** Without the MOve approach, the total expected number of neighbors maintained at node  $p$  is: given by

$$Nbrs_{Worst-Case}(p) = f(N) + \sum_{i=1}^k f(R_i) \quad (1)$$

Now, with the MOve variant we are analyzing, the total expected number of neighbors for a node can be formally represented as union of  $k + 1$  sets as:

$$Nbrs_{MOve}(p) = |Nbrs_N(p) \cup Nbrs_{R_1}(p) \cup Nbrs_{R_2}(p) \cup \dots \cup Nbrs_{R_k}(p)|.$$

This can be written as:

$$\begin{aligned} &= |Nbrs_N(p)| + \sum_{i=1}^k |Nbrs_{R_i}(p)| \\ &\quad - [(\sum_{i=1}^k |Nbrs_N(p) \cap Nbrs_{R_i}(p)|) + (\sum_{i \neq j, 1 \leq i, j \leq k} |Nbrs_{R_i}(p) \cap Nbrs_{R_j}(p)|)] \\ &\quad + [(\sum_{i \neq j, 1 \leq i, j \leq k} |Nbrs_N(p) \cap Nbrs_{R_i}(p) \cap Nbrs_{R_j}(p)|) \\ &\quad \quad + (\sum_{i \neq j \neq l \neq i, 1 \leq i, j, l \leq k} |Nbrs_{R_i}(p) \cap Nbrs_{R_j}(p) \cap Nbrs_{R_l}(p)|)] \\ &\quad - \dots \\ &\quad \pm |Nbrs_N(p) \cup \dots \cup Nbrs_{R_k}(p)| \end{aligned} \quad (2)$$

In order to simplify this, consider an individual term of the type

$|Nbrs_{A_1} \cup \dots \cup Nbrs_{A_m}|$ , where each  $A_j$  is either a unique  $R_i$  or  $N$ . Now consider an arbitrary neighbor  $q$  of  $p$  that is in group  $A_1$ . Consider the event  $E$  that for a given  $j (\neq i)$ , the same neighbor  $q$  (1) also belongs to group  $A_j$  and (2) is a neighbor of  $p$  in group  $A_j$  (i.e., appears in  $Nbrs_{A_j}(p)$ ).

Due to assumptions (2) and (3) in the above theorem, we have that the probability of the above event  $E$  is simply

$$Pr[E] = \frac{|A_j|}{N} \cdot \frac{f(A_j)}{|A_j|} = \frac{f(A_j)}{N}$$

Thus, the individual term of the type of the type  $|Nbrs_{A_1} \cup \dots \cup Nbrs_{A_m}|$  in fact has a value of:

$$|Nbrs_{A_1} \cup \dots \cup Nbrs_{A_m}| = f(A_1) \cdot \prod_{j=2}^m \left(\frac{A_j}{N}\right) = \frac{\prod_{j=1}^m f(A_j)}{N^{m-1}} \quad (3)$$

Substituting equation (3) into equation (2) and using equation (1) above, we get

$$\begin{aligned} Nbrs_{MOve}(p) &= Nbrs_{Worst-Case}(p) \\ &\quad - \frac{1}{N} \cdot [(\sum_{i=1}^k (f_N(p) \cdot f_{R_i}(p)) + (\sum_{i \neq j, 1 \leq i, j \leq k} (f_{R_i}(p) \cdot f_{R_j}(p)))] \\ &\quad + \frac{1}{N^2} \cdot [(\sum_{i \neq j, 1 \leq i, j \leq k} (f_N(p) \cdot f_{R_i}(p) \cdot f_{R_j}(p))) \\ &\quad \quad + (\sum_{i \neq j \neq l \neq i, 1 \leq i, j, l \leq k} (f_{R_i}(p) \cdot f_{R_j}(p) \cdot f_{R_l}(p)))] \\ &\quad - \dots \\ &\quad \pm [f_N(p) \cdot \prod_{i=1}^k (f_{R_i}(p))] \end{aligned}$$

By exchanging the  $Nbrs$  terms, and taking  $f_N(p)$  common on the other side, we simplify to calculate the number of links saved by using MOve as:

$Nbrs_{Worst-Case}(p) - Nbrs_{MOve}(p)$ , which is:

$$\begin{aligned}
&= f_N(p) \cdot \left[ \frac{1}{N} \cdot \sum_{i=1}^k f_{R_i}(p) - \frac{1}{N^2} \cdot \sum_{i \neq j, 1 \leq i, j, k} (f_{R_i}(p) \cdot f_{R_j}(p)) + \dots \pm \frac{1}{N^{k-1}} \cdot \prod_{i=1}^k (f_{R_i}(p)) \right] \\
&\quad + \left[ \frac{1}{N^2} \cdot \sum_{i \neq j, 1 \leq i, j, k} (f_{R_i}(p) \cdot f_{R_j}(p)) - \dots \mp \frac{1}{N^{k-1}} \cdot \prod_{i=1}^k (f_{R_i}(p)) \right] \\
&= f_N(p) \cdot \left[ 1 - \prod_{i=1}^k \left( 1 - \frac{f_{R_i}(p)}{N} \right) \right] + \left[ \prod_{i=1}^k \left( 1 - \frac{f_{R_i}(p)}{N} \right) - 1 + \sum_{i=1}^k \left( \frac{f_{R_i}(p)}{N} \right) \right]
\end{aligned} \tag{4}$$

The above result consists of two terms (each within square braces). The second of these two terms can be shown to be  $\geq 0$  (by using telescoping), and the first term is clearly positive. Finally, the first term is linear in  $f_N(p)$ , as desired. This proves (a) and (b). To prove (c) for equi-sized groups at node  $p$ , substitute  $f_{R_i}(p) = f_R(p)$  for all  $i$  in equation (4) above. Then, we get that  $(Nbrs_{Worst-Case}(p) - Nbrs_{MOve}(p))$  is:

$$\begin{aligned}
&= f_N(p) \cdot \left( 1 - \left( 1 - \frac{f_R(p)}{N} \right)^k \right) + \left( 1 - \frac{f_R(p)}{N} \right)^k - 1 + \frac{k \cdot f_R(p)}{N} \\
&\simeq f_N(p) \cdot \frac{k \cdot f_R(p)}{N} - \frac{k \cdot f_R(p)}{N} + \frac{k \cdot f_R(p)}{N} \\
&= f_N(p) \cdot \frac{k \cdot f_R(p)}{N}
\end{aligned}$$

This proves (c). <<

## 5 Preliminary Evaluation

Our algorithms are implemented in Java, as a discrete event simulation. The GT-ITM [3] random topology generator, following the stub-transit model, is used to provide an underlying internetwork to our simulations. The end-to-end latency of a message corresponds to the shortest path between the sender and receiver nodes.

**Non-Application Link Clustering.** For this experiment we used a topology with 520 nodes, a protocol period for the failure detection mechanism of 1 minute. Parameter  $U$ , which determines the number of protocol periods that a node will allow an unchanged list before randomly refreshing it, has a value of 1. Each node stores a strict maximum of 50 links. Figure 6 shows how the clustering coefficient changes with time. After only 50 minutes, the links among the nodes show a very low degree of clustering. It is not as low as a directed random graph with similar degree, but it is more than enough for our purposes.

**Connectivity.** In this experiment, we try to break the connectivity of the overlay by taking it to an extreme scenario. The basic parameters are the same as in the previous experiment. Figure 7 measures the size of the largest strongly connected component in the overlay. This size is equal to the total number of nodes when the overlay is not partitioned. We vary the number of groups from 0



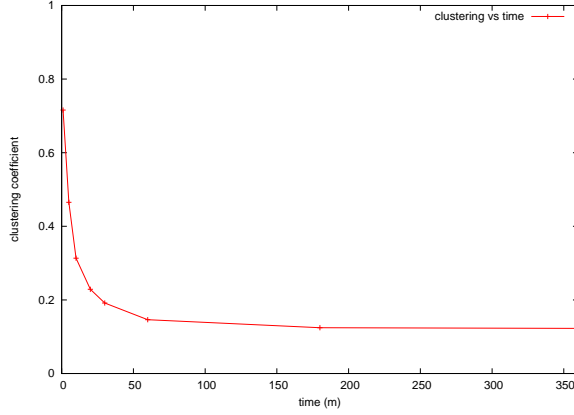


Figure 6: Clustering coefficient vs time.

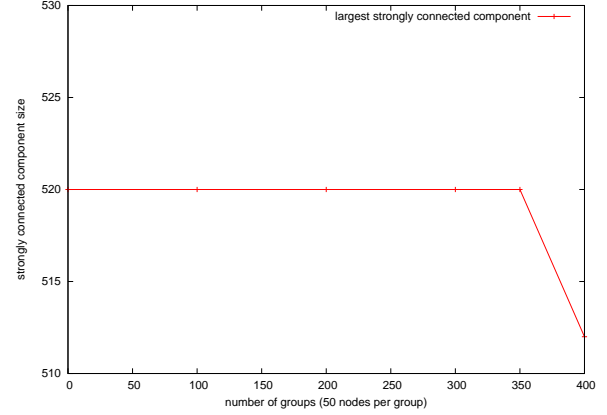


Figure 7: The largest connected component is the overlay itself, until it reaches an overutilization with 400 groups composed of 50 nodes each.

to 400 (each group is composed of 50 nodes, the  $k$  parameter is set to  $\lceil \ln(50) \rceil$ ). As the plot shows, the overlay maintains strong connectivity until the number of groups approaches 400. In this case, we notice a small decrease in size of the largest component, which is due to the overlay partition.

**Subgraphs.** For this experiment, we use a 1000 node network and an application running during 2 hours with one group. We measure the characteristic path length of this group [17]. The characteristic path length is the average of the shortest path over all node pairs. The experiment is run several times varying the subgroup size from 5 to 500 members. The  $k$  parameter (i.e., the number of application links for this group on each node) is set to  $\lceil \ln(\text{groupsize}) \rceil$ . Figure 8 shows that the characteristic path length grows slowly with the group size. Even for 500 nodes it is only 3.27. This shows that the creation of one non-application link at each node of the group, using random walks, achieves its objective of providing a small number of expected hops between any pair of nodes of the group. Note also that characteristic path length follows closely the logarithm with base  $k$ .

**Benefits of Link Sharing.** When the platform contains many groups the probability of intersections grows. In this case, application links can be shared by multiple groups. Figure 9(a) shows the results of simulations upon 520 nodes with 60 groups of 100 nodes each. Parameter  $k$  is set to 5. For readability, only 100 nodes are shown in the figure. It shows that the number of existing application links is much lower than the worst case (which is  $k$  times the number of groups to which a node belongs). This result is due to link-sharing across group intersections, which allows the overlay to use fewer application links when it is solicited by the application. Figure 9(b) repeats the experiment, but this time we correlate groups, by distributing the nodes among the groups using a normal distribution with mean 260 and a standard deviation of 104. This case shows that the number of application links at each node grows at a lower rate than the worst case, thanks to link sharing.

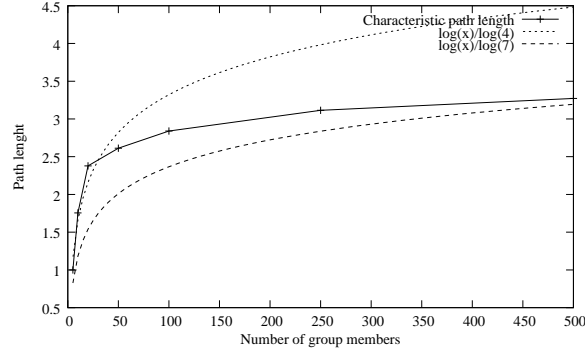


Figure 8: The degree of a group graph is  $\log(\text{groupsize})$  as expected.

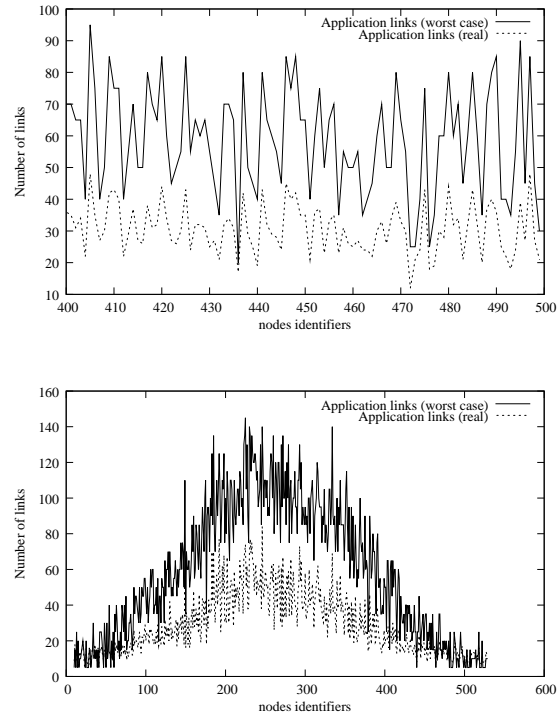


Figure 9: Link sharing among groups. The left plot shows how link sharing is below the worst case when groups are uniformly distributed among nodes. The right plot shows the case when groups are correlated.

**Update Propagation.** Figure 10 shows little variation in update propagation for 100, 200, 300, 400, and 450 replicas. This experiment was done on a 520-node network, with a maximum shortest

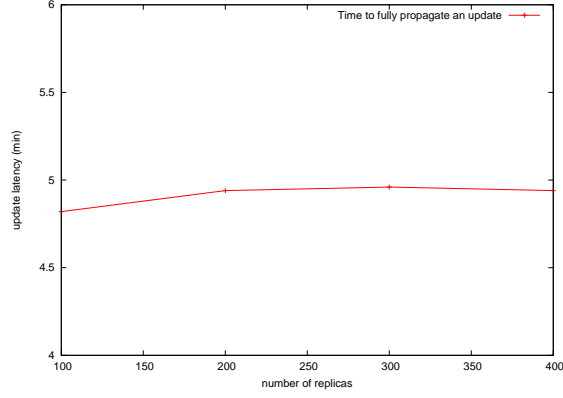


Figure 10: Time taken to propagate an update on all replicas.

path of 290 (ms). In other words, applying the update to all replicas took at most 5 minutes, as the graph shows. The latency stays the same, even though the number of replicas increases, because the characteristic path length for each case is around 3.

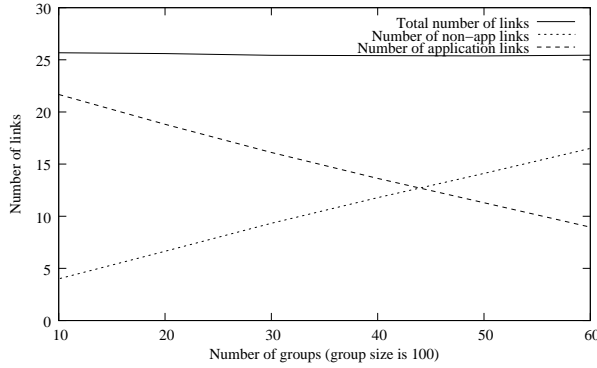


Figure 11: Controlling the clustering

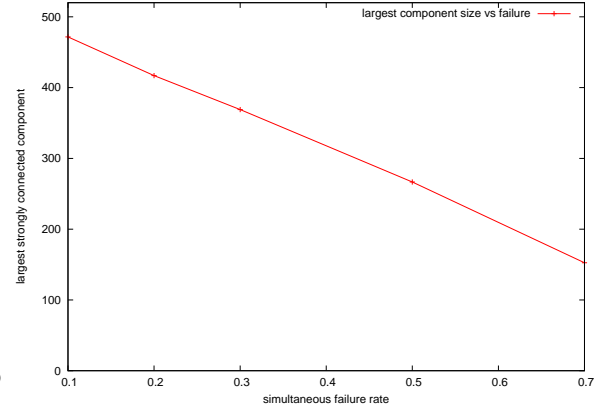


Figure 12: Resilience to simultaneous failure. Overlay partitioning was experienced only when the node failure rate was 0.70

**Twisting the Overlay.** Figure 11 illustrates how the overlay reacts to the application needs (which may differ in the number and size of groups to be created). Simulations were run on a 520-node network, with a varying number of groups having a fixed group size set to 100 (and 5 as  $k$  parameter). The results show that the total number of links is almost constant, while the border between application links and non-application links moves. The creation of groups leads to an increase of the number of application links, which progressively replace the non-application links.

**Resilience to Node Failure.** In this set of experiments, each node may die with a probability of 0.10, 0.20, 0.30, 0.50 and 0.70. The value on the graph is the average over 5 simulations for each probability. Below 0.70, the largest strongly connected component is always the size of the remaining overlay (i.e., the overlay remains strongly connected). With 0.70 death probability, we experienced small partitioning: 2 or 3 nodes were disconnected on some of the runs. Figure 12 shows these results.

## 6 Conclusion

The way the overlay is built is important for *peer-to-peer* applications performance. This paper describes *MOve*, a malleable overlay that applications can “twist”. If an application needs efficient communications between two nodes, it only has to put them in a same group. *MOve* ensures that there are few hops (in the overlay) between two nodes within a same group.

Results showed that the constructed graphs are strongly connected and have a good degree distribution. While an application is running, when groups of nodes are created, the number of application links increases while the number of non-application links decreases, leaving the total number of link on each node almost constant.

We need to further experiment the algorithms presented in this paper. One of the targeted applications is a data-sharing service using replica groups. We plan to implement *MOve* within the *JuxMem* grid data-sharing service [1] using the update mechanism to perform data replication. This will provide the ability to perform extensive evaluations.

## References

- [1] Gabriel Antoniu, Jean-François Deverge, and Sébastien Monnet. How to bring together fault tolerance and data consistency to enable Grid data sharing. *Concurrency and Computation: Practice and Experience*, (17), 2006. To appear.
- [2] Bela Bollobas. *Random Graphs, Second Edition*. Cambridge University Press, United Kingdom, 2001.
- [3] Kenneth L. Calvert, Matthew B. Doar, and Ellen W. Zegura. Modeling Internet topology. *IEEE Communications Magazine*, 35(6):160–163, June 1997.
- [4] Gianni Di Caro and Marco Dorigo. Antnet: Distributed stigmergetic control for communications networks. *Journal of Artificial Intelligence Research*, 9:317–365, 1998.
- [5] Ian Clarke, Oskar Sandberg, Brandon Wiley, and Theodore W. Hong. Freenet: A distributed anonymous information storage and retrieval system. In *Proc. International Workshop on Design Issues in Anonymity and Unobservability*, number 2009, pages 46–66, Berkeley, CA, USA, July 2000.
- [6] Abhinandan Das, Indranil Gupta, and Ashish Motivala. SWIM: Scalable Weakly-consistent Infection-style Process Group Membership Protocol. In *Proc. DSN 02*, Washington DC, June 2002.
- [7] Sidath Handurukande, Anne-Marie Kermarrec, Fabrice Le Fessant, and Laurent Massoulié. Exploiting semantic clustering in the eDonkey p2p network. In *Proc. SIGOPS European Workshop*, pages 109–114, Leuven, Belgium, September 2004.
- [8] P. Hutto and M. Ahamad. Slow memory: Weakening consistency to enhance concurrency in distributed shared memories. In *Proc. ICDCS*, pages 302–311, 1990.

- [9] Mårk Jelasity and Ozalp Babaoglu. T-man: Fast gossip-based construction of large-scale overlay topologies. Technical Report UBLCS-2004-7, University of Bologna, Mura Anteo Zamboni 7 40127 Bologna (Italy), May 2004.
- [10] D. Kostoulas and D. Psaltoulis et al. Decentralized schemes for size estimation in large and dynamic groups. In *Proc. IEEE NCA*, July 2005.
- [11] Andy Oram. *Peer-to-Peer: Harnessing the Power of Disruptive Technologies*, chapter Gnutella, pages 94–122. O'Reilly, May 2001.
- [12] Antony I. T. Rowstron and Peter Druschel. Pastry: Scalable, decentralized object location, and routing for large-scale peer-to-peer systems. In *Proc. Middleware 2001*, Heidelberg, Germany, November 2001.
- [13] Antony I. T. Rowstron, Anne-Marie Kermarrec, Miguel Castro, and Peter Druschel. SCRIBE: The design of a large-scale event notification infrastructure. In *Proc. Networked Group Communication*, pages 30–43, 2001.
- [14] Ion Stoica, Robert Morris, David Karger, Frans Kaashoek, and Hari Balakrishnan. Chord: A scalable peer-to-peer lookup service for Internet applications. In *Proc. SIGCOMM 2001*, pages 149–160, San Diego, CA, August 2001.
- [15] Spyros Voulgaris and Maarten Van Steen Daniela Gavida. Cyclon: Inexpensive membership management for unstructured p2p overlays. *Journal of Network and Systems Management*, 13(2), june 2005.
- [16] Spyros Voulgaris and Maarten van Steen. Epidemic-style management of semantic overlays for content-based searching. In *Proc. 11th Euro-Par*, LNCS, pages 1143–1152, Lisboa, Portugal, August 2005. Springer-Verlag.
- [17] D. J. Watts and S. H. Strogatz. Collective dynamics of "small-world" networks. *Nature*, (393):440–442, June 1998.
- [18] KaZaA. <http://www.kazaa.com/>.



---

Unité de recherche INRIA Lorraine, Technopôle de Nancy-Brabois, Campus scientifique,  
615 rue du Jardin Botanique, BP 101, 54600 VILLERS LÈS NANCY  
Unité de recherche INRIA Rennes, Irista, Campus universitaire de Beaulieu, 35042 RENNES Cedex  
Unité de recherche INRIA Rhône-Alpes, 655, avenue de l'Europe, 38330 MONTBONNOT ST MARTIN  
Unité de recherche INRIA Rocquencourt, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex  
Unité de recherche INRIA Sophia-Antipolis, 2004 route des Lucioles, BP 93, 06902 SOPHIA-ANTIPOLIS Cedex

---

Éditeur  
INRIA, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex (France)  
<http://www.inria.fr>  
ISSN 0249-6399